

The IUCr Diffraction Data Deposition Working Group

John R. Helliwell

Activities

Brian McMahon

John R. Helliwell

School of Chemistry
University of Manchester
Oxford Road
Manchester M13 9PL, UK

john.helliwell@manchester.ac.uk

Activities



Brian McMahon

Research & Development Officer
International Union of Crystallography
5 Abbey Square
Chester CH1 2HU, UK

bm@iucr.org

[http://forums.iucr.org/ viewforum.php?f=7](http://forums.iucr.org/viewforum.php?f=7)

Abstract

It is increasingly important to deposit the raw data from scattering experiments; valuable information gets lost when only structure factors are deposited. Some research centres, e.g. synchrotron and neutron facilities, are fully aware of the need to archive raw data. Diamond Light Source, ISIS, the European PaN (Photon and Neutron facilities) and the Australian TARDIS initiatives are exemplars of good practice. Local university data repositories are also being developed, e.g. at the University of Manchester, relevant to laboratory X-ray diffraction data archiving. A globally registered digital identifier (such as DOI) for each significant dataset can therefore be obtained and linked to any related publication. A Working Group on the development of standards for the representation of data and associated metadata to permit the routine deposition of such raw data has been launched by the IUCr Executive Committee. Its title is 'Diffraction Data Deposition Working Group of the IUCr'. Its members include representatives from IUCr Commissions, synchrotron facilities, academic laboratories, structural databases, and inter-Scientific Union bodies on data and publications. There is a forum for 'Public input on diffraction data deposition' on the IUCr web site at <http://forums.iucr.org>, to which all interested parties are invited to contribute.

Mining the data

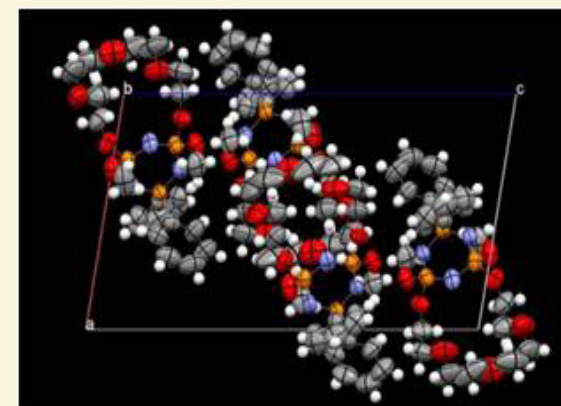
The phrase 'data mining' has often been used to describe the extraction of information from large volumes of available data (e.g. performing a search across all the entries in the Cambridge Structural Database or Protein Data Bank). In fact, it also suggests a valuable metaphor for the processing of experimental data through to a published structural model, especially in terms of volume. There is often a huge amount of raw data from an experiment (several gigabytes of image data), similar to the amount of earth that must be excavated in a mining operation. The processed data are much smaller in volume (perhaps only a few megabytes of structure-factor listing), and from this high-grade 'ore' may be refined or extracted the few kilobytes of 'precious metal' that we can identify with the structural coordinates and displacement parameters of the final model. In most mining operations, once processed, the various spent materials are tossed aside as spoil; yet they may contain other minerals or materials of less immediately obvious value, yet interesting in their own right. In the same way, the raw data that are now routinely discarded may contain information about other scientific features that do not interest the primary investigator. But, in principle, they are there to be unearthed by future prospectors. The question is: can we afford such secondary mining processes?

Why publish data?

Some reasons:

- To enhance the reproducibility of a scientific experiment
- To verify or support the validity of deductions from an experiment
- To safeguard against error
- To safeguard against fraud
- To allow other scholars to conduct further research based on experiments already conducted
- To allow reanalysis at a later date, especially to extract 'new' science as new techniques are developed
- To provide example materials for teaching and learning
- To provide long-term preservation of experimental results and future access to them
- To permit systematic collection for comparative studies

All are – in some measure – ‘good’ reasons. The question is: how much effort, time and money is it reasonable to spend for the potential benefits that might accrue from publishing particular types of data (raw, processed, derived)?



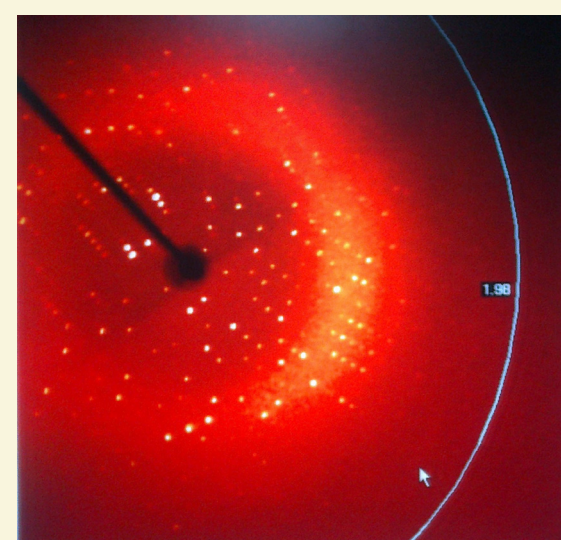
Derived data are those summarizing the refined structural model.

```

        _mipx_ptr0_1600_1280_1024_768_512_256_128_64_32_16_8_4_2_1
        _mipx_ptr0_maxdim 130.0
        _mipx_ptr0_maxdim 120.0
        _mipx_ptr0_maxdim 110.0
        _mipx_ptr0_maxdim 100.0
        _mipx_ptr0_maxdim 90.00
        _mipx_ptr0_maxdim 80.00
        _mipx_ptr0_maxdim 70.00

        _mipx_ptr0_maxdim 60.00
        _mipx_ptr0_maxdim 50.00
        _mipx_ptr0_maxdim 40.00
        _mipx_ptr0_maxdim 30.00
        _mipx_ptr0_maxdim 20.00
        _mipx_ptr0_maxdim 10.00
        _mipx_ptr0_maxdim 5.00
        _mipx_ptr0_maxdim 2.50
        _mipx_ptr0_maxdim 1.25
        _mipx_ptr0_maxdim 0.625
        _mipx_ptr0_maxdim 0.3125
        _mipx_ptr0_maxdim 0.15625
        _mipx_ptr0_maxdim 0.078125
        _mipx_ptr0_maxdim 0.0390625
        _mipx_ptr0_maxdim 0.01953125
        _mipx_ptr0_maxdim 0.009765625
        _mipx_ptr0_maxdim 0.0048828125
        _mipx_ptr0_maxdim 0.00244140625
        _mipx_ptr0_maxdim 0.001220703125
        _mipx_ptr0_maxdim 0.0006103515625
        _mipx_ptr0_maxdim 0.00030517578125
        _mipx_ptr0_maxdim 0.000152587890625
        _mipx_ptr0_maxdim 7.62939453125e-05
        _mipx_ptr0_maxdim 3.814697265625e-05
        _mipx_ptr0_maxdim 1.9073486328125e-05
        _mipx_ptr0_maxdim 9.5367431640625e-06
        _mipx_ptr0_maxdim 4.76837158203125e-06
        _mipx_ptr0_maxdim 2.384185791015625e-06
        _mipx_ptr0_maxdim 1.1920928955078125e-06
        _mipx_ptr0_maxdim 5.9604644775390625e-07
        _mipx_ptr0_maxdim 2.98023223876953125e-07
        _mipx_ptr0_maxdim 1.490116119384765625e-07
        _mipx_ptr0_maxdim 7.450580596923828125e-08
        _mipx_ptr0_maxdim 3.7252902984619140625e-08
        _mipx_ptr0_maxdim 1.86264514923095703125e-08
        _mipx_ptr0_maxdim 9.31322574615478515625e-09
        _mipx_ptr0_maxdim 4.656612873077392578125e-09
        _mipx_ptr0_maxdim 2.3283064365386962890625e-09
        _mipx_ptr0_maxdim 1.16415321826934814453125e-09
        _mipx_ptr0_maxdim 5.82076609134674072265625e-10
        _mipx_ptr0_maxdim 2.910383045673370361328125e-10
        _mipx_ptr0_maxdim 1.4551915228366851806640625e-10
        _mipx_ptr0_maxdim 7.2759576111833259033203125e-11
        _mipx_ptr0_maxdim 3.63797880559166295166015625e-11
        _mipx_ptr0_maxdim 1.818989402795831475830078125e-11
        _mipx_ptr0_maxdim 9.094947013979157379150390625e-12
        _mipx_ptr0_maxdim 4.5474735069895786895751953125e-12
        _mipx_ptr0_maxdim 2.27373675349478934478759765625e-12
        _mipx_ptr0_maxdim 1.136868376747394672393798828125e-12
        _mipx_ptr0_maxdim 5.684341883736973361968994140625e-13
        _mipx_ptr0_maxdim 2.8421709418684866809844970703125e-13
        _mipx_ptr0_maxdim 1.42108547093424334049224853515625e-13
        _mipx_ptr0_maxdim 7.10542735467121670246124267578125e-14
        _mipx_ptr0_maxdim 3.552713677335608351230621337890625e-14
        _mipx_ptr0_maxdim 1.7763568386678041756153106689453125e-14
        _mipx_ptr0_maxdim 8.8817841933390208780765533447265625e-15
        _mipx_ptr0_maxdim 4.44089209666951043903827667236328125e-15
        _mipx_ptr0_maxdim 2.220446048334755219519138336181640625e-15
        _mipx_ptr0_maxdim 1.1102230241673776097595691680908203125e-15
        _mipx_ptr0_maxdim 5.5511151208368880487978458404541015625e-16
        _mipx_ptr0_maxdim 2.77555756041844402439892292022705078125e-16
        _mipx_ptr0_maxdim 1.387778780209222012199461460113525390625e-16
        _mipx_ptr0_maxdim 6.938893901046110060997307300567578125e-17
        _mipx_ptr0_maxdim 3.4694469505230550304986536502837890625e-17
        _mipx_ptr0_maxdim 1.73472347526152751524932682514189453125e-17
        _mipx_ptr0_maxdim 8.67361737630763757624663412570947265625e-18
        _mipx_ptr0_maxdim 4.336808688153818788123317062854736328125e-18
        _mipx_ptr0_maxdim 2.168404344076909394061658531427368125e-18
        _mipx_ptr0_maxdim 1.0842021720384546970308292657136840625e-18
        _mipx_ptr0_maxdim 5.4210108601922734851541463285684203125e-19
        _mipx_ptr0_maxdim 2.71050543009613674257707316428421015625e-19
        _mipx_ptr0_maxdim 1.355252715048068371288536582142105078125e-19
        _mipx_ptr0_maxdim 6.776263575240341856442682910710525390625e-20
        _mipx_ptr0_maxdim 3.388131787620170928221341455355262890625e-20
        _mipx_ptr0_maxdim 1.69406589381008546411067072767763140625e-20
        _mipx_ptr0_maxdim 8.47032946905042732055335363838815703125e-21
        _mipx_ptr0_maxdim 4.235164734525213660276676819194078515625e-21
        _mipx_ptr0_maxdim 2.11758236726260683013833840959703928125e-21
        _mipx_ptr0_maxdim 1.058791183631303415069169204798519640625e-21
        _mipx_ptr0_maxdim 5.293955918156517075345846023992598125e-22
        _mipx_ptr0_maxdim 2.6469779590782585376729230119962990625e-22
        _mipx_ptr0_maxdim 1.32348897953912926883646150599814953125e-22
        _mipx_ptr0_maxdim 6.61744489769564634418230752999074765625e-23
        _mipx_ptr0_maxdim 3.308722448847823172091153764995373828125e-23
        _mipx_ptr0_maxdim 1.65436122442391158604557688249768690625e-23
        _mipx_ptr0_maxdim 8.27180612211955793022788441248843453125e-24
        _mipx_ptr0_maxdim 4.135903061059778965113942206244217265625e-24
        _mipx_ptr0_maxdim 2.0679515305298894825569711031221086328125e-24
        _mipx_ptr0_maxdim 1.03397576526494474127848555156105431640625e-24
        _mipx_ptr0_maxdim 5.169878826324723706392427757805271875e-25
        _mipx_ptr0_maxdim 2.5849394131623618531962138789026359375e-25
        _mipx_ptr0_maxdim 1.2924697065811809265981069394513178125e-25
        _mipx_ptr0_maxdim 6.4623485329059046329905346972565890625e-26
        _mipx_ptr0_maxdim 3.23117426645295231649526734862829453125e-26
        _mipx_ptr0_maxdim 1.615587133226476158247633674314147265625e-26
        _mip
```

Processed data are the structure factors or Rietveld profiles reduced by standard techniques from the raw experimental output.



Raw data are those generated by the experimental apparatus.

IUCr journal policy

(1) Derived data

For crystal/molecular structures with small unit cell

- Atomic coordinates, anisotropic displacement parameters, molecular geometry and intermolecular contacts
 - Experimental parameters, unit-cell dimensions, space group information
 - Reference and modulated structure subsystems for aperiodic composite structures
- must be supplied in CIF format as an integral part of article submission and are freely available for download.

For biological macromolecular structures

- Atomic coordinates, anisotropic or isotropic displacement parameters, space group information, secondary structure and information about biological functionality must be deposited with the Protein Data Bank before or in concert with article publication; the article will link to the PDB deposition using the PDB reference code.
- Relevant experimental parameters, unit-cell dimensions are required as an integral part of article submission and are published within the article.

(2) Processed experimental data

For crystal/molecular structures with small unit cell

- Structure factors
- Rietveld profiles

must be supplied in CIF format as an integral part of article submission and are freely available for download. *SHELXL* instruction files are also required for validation.

For biological macromolecular structures

- Structure factors

must be deposited with the Protein Data Bank before or in concert with article publication; the article will link to the PDB deposition using the PDB reference code.

(3) Primary experimental data

For small-unit-cell crystal/molecular structures and macromolecular structures

IUCr journals have *no* current binding policy regarding publication of diffraction images or similar raw data entities.

However, the journals welcome efforts made to preserve and provide primary experimental data sets, and encouragement is often provided in Notes for Authors. For example, submission guidelines for *Acta Crystallographica Section D* state:

Fibre data should contain appropriate information such as a photograph of the data. As primary diffraction data cannot be satisfactorily extracted from such figures, the basic digital diffraction data should be deposited.

Authors are encouraged to make arrangements for the diffraction data images for their structure to be archived and available on request.

For articles that present the results of powder diffraction profile fitting or refinement (Rietveld) methods, the primary diffraction data, i.e. the numerical intensity of each measured point on the profile as a function of scattering angle, should be deposited.

Furthermore, many IUCr Commissions are interested in the possibility of establishing community practices for the orderly retention and referencing of such data sets, and the IUCr would like to see such data sets become part of the routine record of scientific research in the future, to the extent that this proves feasible and cost-effective.

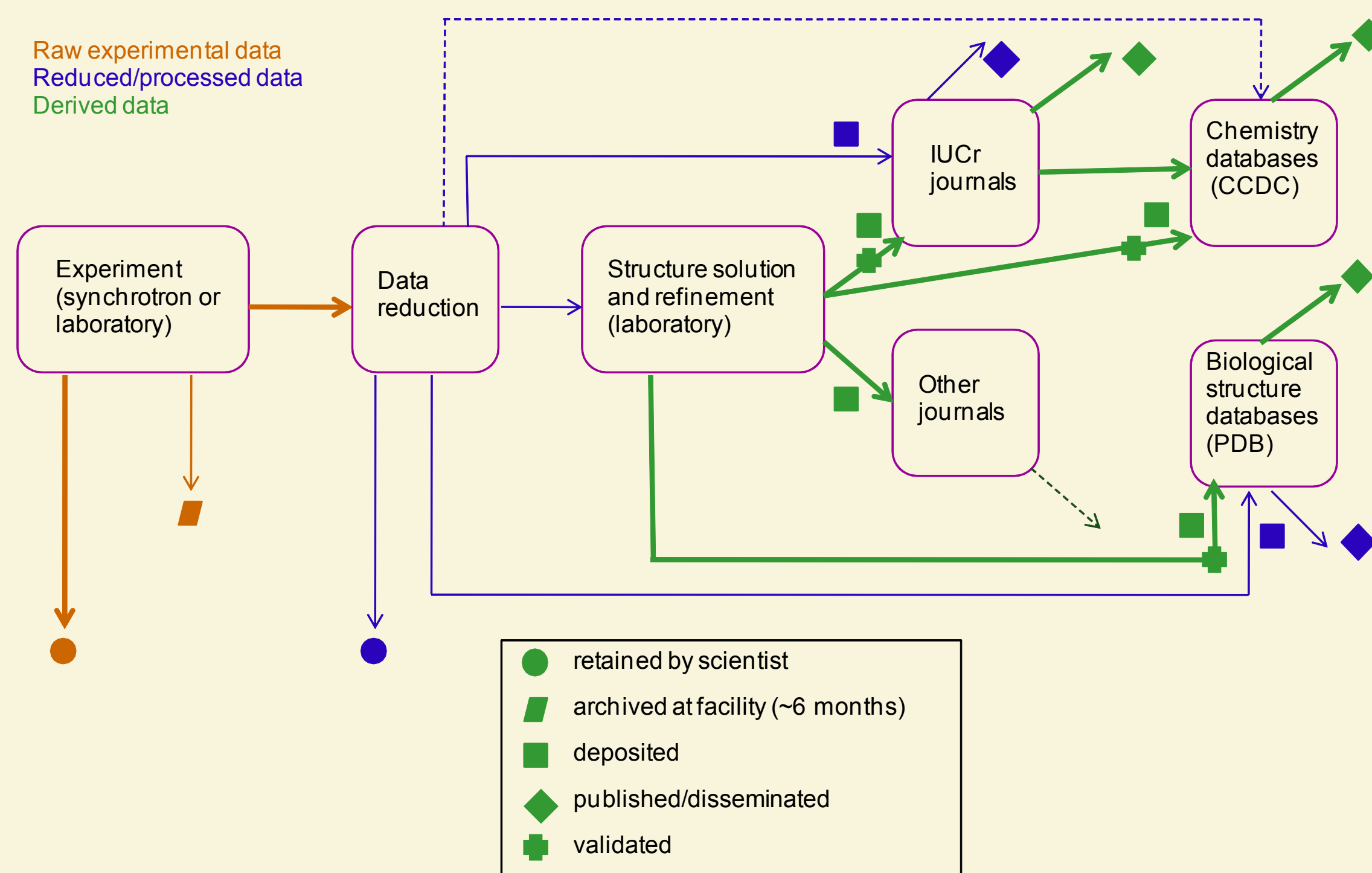
Policy of funding bodies

Increasingly, funding agencies are requesting or requiring data management policies (including provision for retention and access) to be taken into account when awarding grants. See e.g. the Research Councils UK Common Principles on Data Policy (<http://www.rcuk.ac.uk/research/Pages/DataPolicy.aspx>) and the Digital Curation Centre overview of funding policies in the UK (<http://www.dcc.ac.uk/resources/policy-and-legal/overview-funders-data-policies>).

See also <http://forums.iucr.org/viewtopic.php?f=21&t=58> for discussion on policies relevant to crystallography in other countries.

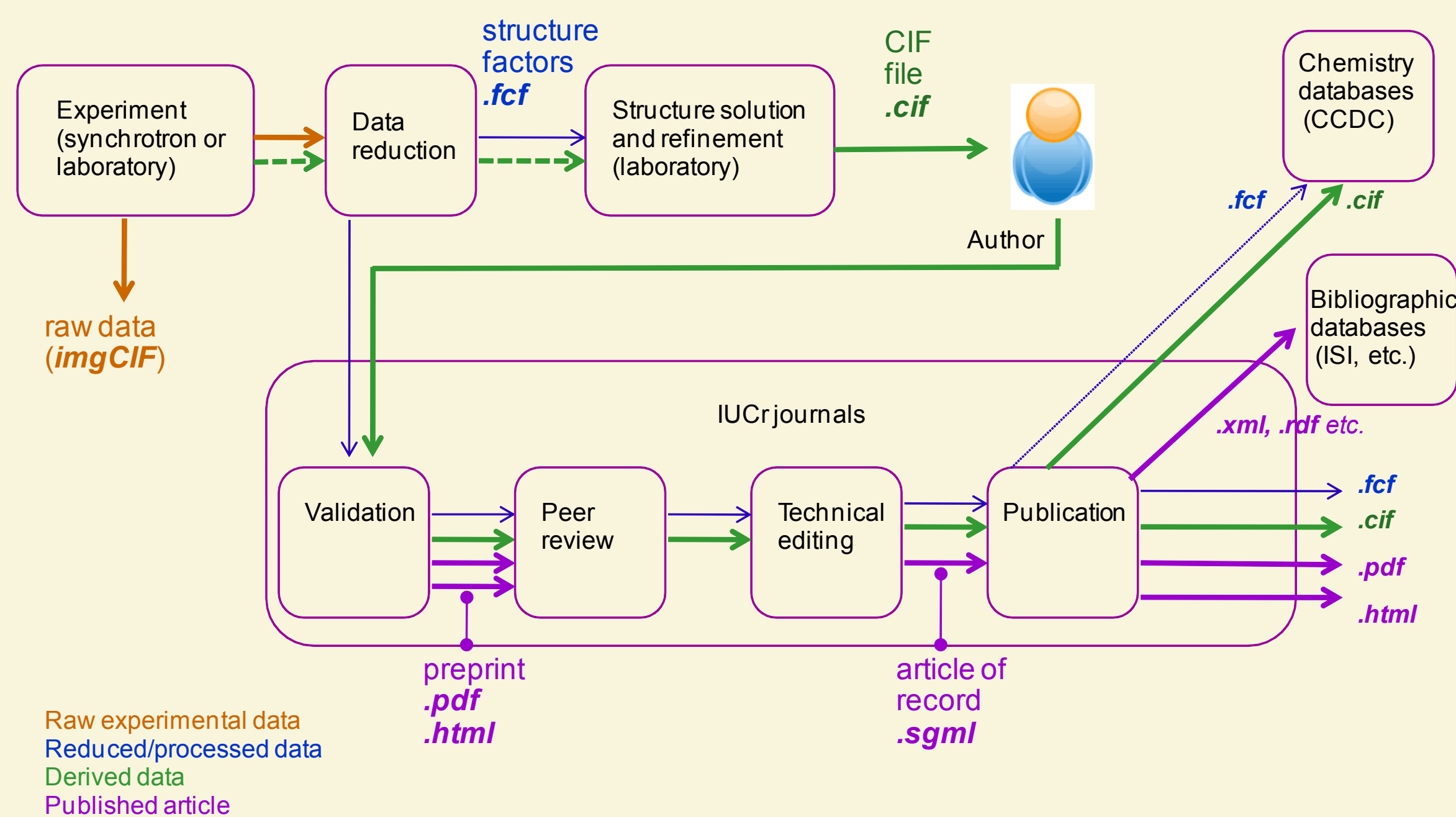
Data flow in crystallography

This schematic indicates the flow of data from experiment, through analysis, publication and deposition of structural results, for a typical single-crystal X-ray structure determination. The relationship between journals and curated databases differs between small-molecule and macromolecular communities. But there are broadly parallel flows, and owing to standard Web-based information access and delivery systems, publishers, databases, laboratories and large-scale facilities all have the potential of working together to create a distributed information ecosystem.



Publication flow in IUCr journals

This schematic shows the information flow associated with publication of crystal structure reports in IUCr journals. The data associated with the publication are handled as outlined in the upper figure, but here we emphasise the role of the author, who interpolates the textual content of a paper directly into a CIF (in the case of *Acta Crystallographica Sections C* and *E*), the value added to the publication through peer review (including the semi-automated analysis and validation performed by the *checkCIF* service), and the multiple format translations through which the data and associated comment comprising an article must pass during the technical editing and publishing processes. It is important that many of the format interconversions can occur losslessly (*i.e.* without decay of information content). In practice, subsets of the information are filtered out during various stages of processing. The article is usually delivered as a HTML or PDF document, with little numerical data retained (except to the extent that it is relevant to the scientific discussion), and less semantic markup than in the input files (thus, PDF is largely an end-stage presentational medium rather than a working format), while the numerical data in the input CIF and structure factor files are usually served as supplementary files.



Working group members

- Steve Androulakis *Representative of TARDIS (Australian repositories for diffraction images)*
- Sol Gruner *Diffuse scattering specialist and Synchrotron Radiation Facility Director*
- John R. Helliwell, **Chair** *IUCr Representative to ICSTI; Chair, IUCr Commission on Journals 1996-2005*
- Loes Kroon-Battenburg *Data processing software developer and user*
- Brian McMahon *IUCr Representative to CODATA*
- Tom Terwilliger *Representative of IUCr Commission on Biological Crystallography*
- John Westbrook *Representative of wwPDB (Worldwide Protein Data Bank)*
- Hans-Josef Weyer *Synchrotron Radiation and Neutron Facility user*

Consultants

- Alun Ashton *Diamond Light Source*
- Herbert Bernstein *Head, imgCIF Dictionary Maintenance Group and member of COMCIFS*
- Frances Bernstein *Observer on data deposition policies*
- Gerard Bricogne *Active software and methods developer*
- Bernhard Rupp *Macromolecular crystallographer*

References and further reading

Research data management policy: UK Research Data Service, *The data imperative: Managing the UK's research data for future use* (2009). Available from:

Research data management policy: e-IRG Data Management Task Force. *Report on Data Management* (2009). Available from:

Research data management policy: High Level Expert Group on Scientific Data *Riding the wave - How Europe can gain from the rising tide of scientific data* (2010). Available from: http://ec.europa.eu/information_society/newsroom/cf/itemlongdetail.cfm?item_id=6204

Domain repository solutions: S. Androulakis, J. Schmidberger, M. A. Bate, R. DeGori, A. Beitz, C. Keong, B. Cameron, S. McGowan, C. J. Porter, A. Harrison, J. Hunter, J. L. Martin, B. Kobe, R. C. J. Dobson, M. W. Parker, J. C. Whisstock, J. Gray, A. Trelor, D. Groenewegen, N. Dickson and A. M. Buckle *Federated repositories of X-ray diffraction images* (2008). *Acta Cryst.* **D64**, 810-814
| doi:10.1107/S090744440908015540 |

Infrastructure architecture: Photon and Neutron Open Data Infrastructure (2012) <http://pan-data.eu>

Persistent identifiers for data: DataCite (2009). <http://www.datacite.org/>